



# Efficient certification of numeric solutions to eigenproblems

Joris van Der Hoeven, Bernard Mourrain

## ► To cite this version:

Joris van Der Hoeven, Bernard Mourrain. Efficient certification of numeric solutions to eigenproblems. MACIS 2017 - 7th International Conference on Mathematical Aspects of Computer and Information Sciences, Nov 2017, Vienna, Austria. pp.81–94, 10.1007/978-3-319-72453-9\_6 . hal-01579079

**HAL Id: hal-01579079**

**<https://hal.science/hal-01579079>**

Submitted on 30 Aug 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Efficient certification of numeric solutions to eigenproblems

JORIS VAN DER HOEVEN

Laboratoire d'informatique, UMR 7161 CNRS  
Campus de l'École polytechnique  
1, rue Honoré d'Estienne d'Orves  
Bâtiment Alan Turing, CS35003  
91120 Palaiseau  
France  
Email: vdhoeven@lix.polytechnique.fr

BERNARD MOURRAIN

Inria Sophia Antipolis Méditerranée  
AROMATH  
2004 route des Lucioles  
06902 Sophia Antipolis  
France  
Email: Bernard.Mourrain@inria.fr

August 30, 2017

---

In this paper, we present an efficient algorithm for the certification of numeric solutions to eigenproblems. The algorithm relies on a mixture of ball arithmetic, a suitable Newton iteration, and clustering of eigenvalues that are close.

KEYWORDS: Ball arithmetic, interval arithmetic, reliable computing, computable analysis

A.M.S. SUBJECT CLASSIFICATION: 65G20, 03F60, 65F99, 37-04

---

## 1. INTRODUCTION

Let  $\mathbb{F}$  be the set of floating point numbers for a fixed precision and a fixed exponent range. We will denote  $\mathbb{F}^{\geq} = \{x \in \mathbb{F} : x \geq 0\}$ . Consider an  $n \times n$  matrix  $M \in \mathbb{F}[i]^{n \times n}$  with complex floating entries. The numeric *eigenproblem* associated to  $M$  is to compute a transformation matrix  $T \in \mathbb{F}[i]^{n \times n}$  and a diagonal matrix  $D \in \mathbb{F}[i]^{n \times n}$  such that

$$D \approx T^{-1}MT. \quad (1)$$

The entries of  $D$  are the approximate eigenvalues and the columns of  $T$  are the approximate eigenvectors of  $M$ . In addition, we might require that  $T$  is normalized. For instance, each of the columns might have unit norm. Alternatively, the norm of the  $i$ -th column may be required to be the same as the norm of the  $i$ -th row of  $T^{-1}$ , for each  $i$ . There are several well-known algorithms for solving the numeric eigenproblem [4].

Unfortunately, (1) is only an approximate equality. It is sometimes important to have rigorous bounds for the distance between the approximate eigenvalues and/or eigenvectors and the genuine ones. More precisely, we may ask for a diagonal matrix  $D_r \in (\mathbb{F}^{\geq})^{n \times n}$  and a matrix  $T_r \in (\mathbb{F}^{\geq})^{n \times n}$  of radii such that there exists a matrix  $T' \in \mathbb{C}^{n \times n}$  for which

$$D' = (T')^{-1}MT$$

is diagonal and

$$\begin{aligned} |D'_{i,i} - D_{i,i}| &\leq (D_r)_{i,i} \\ |T'_{i,j} - T_{i,j}| &\leq (T_r)_{i,j} \end{aligned}$$

for all  $i, j$ . This task will be called the *certification problem* of the numeric solution  $(D, T)$  to the eigenproblem for  $M$ .

It will be convenient to rely on *ball arithmetic* [6, 8], which is a systematic technique for bound computations. When computing with complex numbers, ball arithmetic is more accurate than more classical interval arithmetic [11, 1, 12, 7, 9, 15], especially in multiple precision contexts. We will write  $\mathbb{B} = \mathcal{B}(\mathbb{F}[\mathbf{i}], \mathbb{F}^{\geq})$  for the set of balls  $\mathbf{z} = \mathcal{B}(z_c, z_r) = \{z \in \mathbb{C} : |z - z_c| \leq z_r\}$  with centers  $z_c$  in  $\mathbb{F}[\mathbf{i}]$  and radii  $z_r$  in  $\mathbb{F}^{\geq}$ . In a similar way, we may consider matricial balls  $\mathbf{M} = \mathcal{B}(M_c, M_r) \in \mathcal{B}(\mathbb{F}[\mathbf{i}]^{n \times n}, (\mathbb{F}^{\geq})^{n \times n})$ : given a center matrix  $M_c \in \mathbb{F}[\mathbf{i}]^{n \times n}$  and a radius matrix  $M_r \in (\mathbb{F}^{\geq})^{n \times n}$ , we have

$$\mathbf{M} = \mathcal{B}(M_c, M_r) = \{M \in \mathbb{C}^{n \times n} : \forall i, j, |(M_c)_{i,j} - M_{i,j}| \leq (M_r)_{i,j}\}.$$

Alternatively, we may regard  $\mathcal{B}(M_c, M_r)$  as the set of matrices in  $\mathbb{B}^{n \times n}$  with ball coefficients:

$$\mathcal{B}(M_c, M_r)_{i,j} = \mathcal{B}((M_c)_{i,j}, (M_r)_{i,j}).$$

Standard arithmetic operations on balls are carried out in a reliable way. For instance, if  $\mathbf{u}, \mathbf{v} \in \mathbb{B}$ , then the computation of the product  $\mathbf{w} = \mathbf{u} \mathbf{v}$  using ball arithmetic has the property that  $u v \in \mathbf{w}$  for any  $u \in \mathbf{u}$  and  $v \in \mathbf{v}$ . Given a ball  $\mathbf{z} \in \mathbb{B}$ , it will finally be convenient to write  $\lfloor \mathbf{z} \rfloor \in \mathbb{F}^{\geq}$  and  $\lceil \mathbf{z} \rceil \in \mathbb{F}^{\geq}$  for certified lower and upper bounds of  $|\mathbf{z}|$  in  $\mathbb{F}^{\geq}$ .

In the language of ball arithmetic, it is natural to allow for small errors in the input and replace the numeric input  $M \in \mathbb{F}[\mathbf{i}]^{n \times n}$  by a ball input  $\mathcal{B}(M_c, M_r) \in \mathbb{B}^{n \times n}$ . Then we may still compute a numeric solution

$$D_c \approx T_c^{-1} M_c T_c, \tag{2}$$

for the eigenproblem associated to the center  $M_c$ . Assume that the matrices in  $\mathcal{B}(M_c, M_r)$  are all diagonalizable. The generalized *certification problem* now consists of the computation of a diagonal matrix  $D_r \in (\mathbb{F}^{\geq})^{n \times n}$  and a matrix  $T_r \in \mathbb{F}[\mathbf{i}]^{n \times n}$  such that, for every  $M \in \mathcal{B}(M_c, M_r)$ , there exist  $D \in \mathcal{B}(D_c, D_r)$  and  $T \in \mathcal{B}(T_c, T_r)$  with

$$D = T^{-1} M T.$$

In absence of multiple eigenvalues, known algorithms for solving this problem such as [17, 14] proceed by the individual certification of each eigenvector, which results in an  $O(n^4)$  running time.

Extensions to a cluster of eigenvalues and the corresponding eigenvectors have been considered in [3, 16], with similar  $O(n^4)$  complexity bounds. Fixed points theorem based on interval arithmetic are used to prove the existence of a matrix with a given Jordan block in the matrix interval domain. Such an approach has been exploited for the analysis of multiple roots in [5, 13]. A test that provides an enclosing of all the eigenvalues has been proposed in [10]. Its certification relies on interval and ball arithmetics. The complexity of the test is in  $O(n^3)$  but no iteration converging to the solution of the eigenproblem is described.

In this paper, we present a new algorithm of time complexity  $O(n^3)$  for certifying and enclosing clusters of eigenvectors and eigenvalues in a single step. We also provide an iterative procedure that converges geometrically to clusters of solutions. This convergence is quadratic in the case of single eigenvalues. Our algorithm extends a previous algorithm from [6] to the case of multiple eigenvalues. This yields an efficient test for *approximate eigenvalues* in the sense of the  $\alpha$ -theory [2].

We recall that it is very unlikely that the numeric matrix  $M_c \in \mathbb{F}[i]^{n \times n}$  with complex floating point coefficients has multiple eigenvalues. Indeed, small perturbations of matrices with multiple eigenvalues, as induced by rounding errors, generically only have simple eigenvalues. Consequently, we may assume without loss of generality that the numeric eigenproblem (2) has a reasonably accurate solution (if necessary, we may slightly perturb  $M_c$  and increase  $M_r$  accordingly). Using ball arithmetic, it is straightforward to compute the matricial ball

$$\mathcal{B}(N_c, N_r) = \mathcal{B}(T_c, 0)^{-1} \mathcal{B}(M_c, M_r) \mathcal{B}(T_c, 0).$$

If our numerical algorithm is accurate, then the non diagonal entries of  $\mathcal{B}(N_c, N_r)$  tend to be small, whence  $\mathcal{B}(N_c, N_r)$  can be considered as a small perturbation of a diagonal matrix. If we can estimate how far eigenvalues and eigenvectors of diagonal matrices can drift away under small perturbations, we thus obtain a solution to the original certification problem.

Section 2 introduces notations. In Section 3, we perform a detailed study of the eigenproblem for small perturbations  $M$  of diagonal matrices. We exhibit a Newton iteration for finding the solutions. This iteration has quadratic convergence in the absence of multiple eigenvalues and is also an efficient tool for doubling the precision of a solution. However, in the case of multiple eigenvalues, the eigenproblem is ill-posed. Indeed, by a well-known observation, *any* vector occurs as the eigenvector of a small perturbation of the  $2 \times 2$  identity matrix. The best we can hope for is to group eigenvectors with close eigenvalues together in “clusters” (see also [16]) and only require  $T^{-1} M T$  to be block diagonal. For this reason, we present our Newton iteration in a sufficiently general setting which encompasses block matrices. We will show that the iteration still admits geometric convergence for sufficiently small perturbations and that the blockwise certification is still sufficient for the computation of rigorous error bounds for the eigenvalues. In Section 4, we will present explicit algorithms for clustering and the overall certification problem.

## 2. NOTATIONS

### 2.1. Matrix norms

Throughout this paper, we will use the max norm for vectors and the corresponding matrix norm. More precisely, given a vector  $v \in \mathbb{C}^n$  and an  $n \times n$  matrix  $M \in \mathbb{C}^{n \times n}$ , we set

$$\begin{aligned} \|v\| &= \max \{|v_1|, \dots, |v_n|\} \\ \|M\| &= \max_{\|v\|=1} \|Mv\|. \end{aligned}$$

For a second matrix  $N \in \mathbb{C}^{n \times n}$ , we clearly have

$$\begin{aligned} \|M + N\| &\leq \|M\| + \|N\| \\ \|MN\| &\leq \|M\| \|N\|. \end{aligned}$$

Explicit machine computation of the matrix norm is easy using the formula

$$\|M\| = \max \{|M_{i,1}| + \dots + |M_{i,n}| : 1 \leq i \leq n\}. \quad (3)$$

In particular, when changing certain entries of a matrix  $M$  to zero, its matrix norm  $\|M\|$  can only decrease.

### 2.2. Clustering

Assume that we are given a partition

$$\{1, \dots, n\} = I_1 \amalg \dots \amalg I_p. \quad (4)$$

Such a partition will also be called a *clustering* and denoted by  $I$ . Two indices  $i, j$  are said to belong to the same *cluster* if there exists a  $k$  with  $\{i, j\} \subseteq I_k$  and we will write  $i \sim j$ . Two entries  $M_{i,j}$  and  $M_{i',j'}$  of a matrix  $M \in \mathbb{C}^{n \times n}$  are said to belong to the same *block* if  $i \sim j$  and  $i' \sim j'$ . We thus regard  $M$  as a generalized block matrix, for which the rows and columns of the blocks are not necessarily contiguous inside  $M$ .

A matrix  $M \in \mathbb{C}^{n \times n}$  is said to be *block diagonal* (relative to the clustering) if  $M_{i,j} = 0$  whenever  $i \not\sim j$ . Similarly, we say that  $M$  is *off block diagonal* if  $M_{i,j} = 0$  whenever  $i \sim j$ . For a general  $M \in \mathbb{C}^{n \times n}$ , we define its block diagonal and off block diagonal projections  $\Delta(M) = \Delta^I(M)$  and  $\Omega(M) = \Omega^I(M)$  by

$$\Delta(M)_{i,j} = \begin{cases} M_{i,j} & \text{if } i \sim j \\ 0 & \text{otherwise} \end{cases} \quad \Omega(M)_{i,j} = \begin{cases} 0 & \text{if } i \sim j \\ M_{i,j} & \text{otherwise} \end{cases}$$

By our observation at the end of section 2.1, we have

$$\begin{aligned} \|\Delta(M)\| &\leq \|M\| \\ \|\Omega(M)\| &\leq \|M\|. \end{aligned}$$

For the *trivial clustering*  $I_k = \{k\}$ , the matrices  $\Delta(M)$  and  $\Omega(M)$  are simply the diagonal and off diagonal projections of  $M$ . In that case we will also write  $\Delta^* = \Delta$  and  $\Omega^* = \Omega$ .

### 2.3. Diagonal matrices

Below, we will study eigenproblems for perturbations of a given diagonal matrix

$$D = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}. \quad (5)$$

It follows from (3) that the matrix norm  $\mu = \|D\|$  of a diagonal matrix  $D$  is given by

$$\mu = \max\{|\lambda_1|, \dots, |\lambda_n|\}.$$

It will also be useful to define the *separation number*  $\sigma^* = \sigma^*(D)$  by

$$\sigma^* = \min\{|\lambda_i - \lambda_j| : i \neq j\}.$$

More generally, given a clustering as in the previous subsection, we also define the *block separation number*  $\sigma = \sigma(D) = \sigma^I(D)$  by

$$\sigma = \min\{|\lambda_i - \lambda_j| : i \not\sim j\}$$

This number  $\sigma$  remains high if the clustering is chosen in such a way that the indices  $i, j$  of any two “close” eigenvalues  $\lambda_i$  and  $\lambda_j$  belong to the same cluster. In particular, if  $\sigma > 0$ , then  $\lambda_i = \lambda_j$  implies  $i \sim j$ .

## 3. EIGENPROBLEMS FOR PERTURBED DIAGONAL MATRICES

### 3.1. The linearized equation

Let  $D$  be a diagonal matrix (5). Given a small perturbation

$$M = D + H$$

of  $D$ , where  $H$  is an off diagonal matrix, the aim of this section is to find a small matrix  $E \in \mathbb{C}^{n \times n}$  for which

$$M' = (1 + E)^{-1} M (1 + E)$$

is block diagonal. In other words, we need to solve the equation

$$\Omega((1 + E)^{-1} (D + H) (1 + E)) = 0.$$

When linearizing this equation in  $E$  and  $H$ , we obtain

$$\Omega([D, E] + H) = 0.$$

If  $E$  is strongly off diagonal, then so is  $[D, E]$ , and the equation further reduces to

$$[D, E] = -\Omega(H).$$

This equation can be solved using the following lemma:

LEMMA 1. *Given a matrix  $A \in \mathbb{C}^{n \times n}$  and a diagonal matrix  $D$  with entries  $\lambda_1, \dots, \lambda_n$ , let  $B = \Phi(D, A) \in \mathbb{C}^{n \times n}$  be the strongly off diagonal matrix with*

$$B_{i,j} = \begin{cases} 0 & \text{if } i \sim j \\ \frac{A_{i,j}}{\lambda_j - \lambda_i} & \text{otherwise} \end{cases}$$

*Then  $\|B\| \leq \sigma^{-1} \|A\|$  and*

$$[D, B] = -\Omega(A). \tag{6}$$

**Proof.** The inequality follows from (3) and the definition of  $\sigma$ . One may check (6) using a straightforward computation.  $\square$

### 3.2. The fundamental iteration

In view of the lemma, we now consider the iteration

$$(D, H) \longmapsto (D', H'),$$

where

$$\begin{aligned} E &= \Phi(D, H) \\ M' &= (1 + E)^{-1} (D + H) (1 + E) \\ D' &= \Delta^*(M') \\ H' &= \Omega^*(M') \end{aligned}$$

In order to study the convergence of this iteration, we introduce the quantities

$$\begin{aligned} \mu &= \|D\| & \mu' &= \|D'\| \\ \sigma &= \sigma(D) & \sigma' &= \sigma(D') \\ \eta_1 &= \|\Delta(H)\| & \eta'_1 &= \|\Delta(H')\| \\ \eta_2 &= \|\Omega(H)\| & \eta'_2 &= \|\Omega(H')\| \end{aligned}$$

$$\alpha = \min \left\{ \frac{\sigma}{6\mu}, \frac{1}{4} \right\}.$$

LEMMA 2. For  $\delta \in (0, 1]$ , assume that

$$\begin{aligned}\eta_1 + \eta_2 &\leq \alpha \delta \mu \\ \eta_2 &\leq \alpha \delta \sigma.\end{aligned}$$

Then  $\|D' - D\| \leq \delta \eta_2$  and

$$\begin{aligned}\mu' &\leq \mu + \delta \eta_2 \\ \sigma' &\geq \sigma - 2 \delta \eta_2 \\ \eta_1' &\leq \eta_1 + \delta \eta_2 \\ \eta_2' &\leq \delta \eta_2.\end{aligned}$$

**Proof.** We have

$$\begin{aligned}M' - D &= H + [D, E] + R \\ &= \Delta(H) + R,\end{aligned}$$

where

$$R = E^2(1 + E)^{-1}(D + H)(1 + E) - E(D + H)E + [H, E].$$

Setting  $\varepsilon = \|E\| \leq \sigma^{-1} \eta_2 \leq \alpha \delta \leq \frac{1}{4}$ , the remainder  $R$  is bounded by

$$\begin{aligned}\|R\| &\leq \varepsilon^2 \frac{1}{1 - \varepsilon} (1 + \alpha \delta) \mu (1 + \varepsilon) + \varepsilon (1 + \alpha \delta) \mu \varepsilon + 2(\eta_1 + \eta_2) \varepsilon \\ &= \frac{2\varepsilon^2}{1 - \varepsilon} (1 + \alpha \delta) \mu + 2(\eta_1 + \eta_2) \varepsilon \\ &\leq (4\varepsilon \mu + 2\alpha \delta \mu) \varepsilon \\ &\leq 6\alpha \delta \mu \sigma^{-1} \eta_2 \\ &\leq \delta \eta_2.\end{aligned}$$

Consequently,

$$\begin{aligned}\|D' - D\| &= \|\Delta^*(M' - D)\| = \|\Delta^*(R)\| \\ &\leq \|R\| \leq \delta \eta_2 \\ \eta_1' &= \|\Delta(H')\| = \|\Omega^*(\Delta(M'))\| = \|\Omega^*(\Delta(H + R))\| \\ &\leq \|H + R\| \leq \eta_1 + \delta \eta_2 \\ \eta_2' &= \|\Omega(H')\| = \|\Omega(M')\| = \|\Omega(R)\| \\ &\leq \delta \eta_2.\end{aligned}$$

The inequalities  $\mu' \leq \mu + \delta \eta_2$  and  $\sigma' \geq \sigma - 2 \delta \eta_2$  follow from  $\|D' - D\| \leq \delta \eta_2$ .  $\square$

### 3.3. Convergence of the fundamental iteration

THEOREM 3. Assume that

$$\begin{aligned}\eta_1 + \eta_2 &\leq \frac{1}{8} \alpha \mu \\ \eta_2 &\leq \frac{1}{8} \alpha \sigma.\end{aligned}$$

Then the sequence

$$(D, H), (D', H'), (D'', H''), \dots$$

converges geometrically to a limit  $(M^{(\infty)}, H^{(\infty)})$  with  $\|D^{(\infty)} - M\| \leq \eta_2$  and  $\|H^{(\infty)}\| \leq \eta_1 + \eta_2$ . The matrix  $D^{(\infty)} + H^{(\infty)}$  is block diagonal and there exists a matrix  $\hat{E}$  with  $\|\hat{E}\| \leq 3\sigma^{-1}\eta_2$ , such that

$$D^{(\infty)} + H^{(\infty)} = (1 + \hat{E})^{-1} (D + H) (1 + \hat{E}).$$

**Proof.** Let  $(D^{(i)}, H^{(i)})$  stand for the  $i$ -th fundamental iterate of  $(D, H)$  and  $E^{(i)} = \Phi(H^{(i)}, D^{(i)})$ . Denote  $\mu^{(i)} = \|D^{(i)}\|$ ,  $\sigma^{(i)} = \sigma(D^{(i)})$ ,  $\eta_1^{(i)} = \|\Delta(H^{(i)})\|$  and  $\eta_2^{(i)} = \|\Omega(H^{(i)})\|$ . Let us show by induction over  $i$  that

$$\begin{aligned} \|D^{(i)} - D\| &\leq (1 - \frac{1}{2^i}) \eta_2 \\ \mu^{(i)} &\leq \mu + (1 - \frac{1}{2^i}) \eta_2 \\ \sigma^{(i)} &\geq \frac{1}{2} (1 + \frac{1}{2^i}) \sigma \\ \eta_1^{(i)} &\leq \eta_1 + (1 - \frac{1}{2^i}) \eta_2 \\ \eta_2^{(i)} &\leq \frac{1}{2^i} \eta_2. \end{aligned}$$

This is clear for  $i=0$ . Assume that the induction hypothesis holds for a given  $i$  and let

$$\alpha^{(i)} = \min \left\{ \frac{\sigma^{(i)}}{6\mu^{(i)}}, \frac{1}{4} \right\}$$

Since  $(1 - \frac{1}{2^i}) \eta_2 \leq \frac{1}{32} \mu$ , the induction hypothesis implies

$$\begin{aligned} \mu^{(i)} &\leq 2\mu \\ \sigma^{(i)} &\geq \frac{1}{2} \sigma \\ \alpha^{(i)} &\geq \frac{1}{4} \alpha. \end{aligned}$$

Applying Lemma 2 for  $(D^{(i)}, H^{(i)})$  and  $\delta = \frac{1}{2}$ , we thus find

$$\begin{aligned} \|D^{(i+1)} - D\| &\leq \|D^{(i)} - D\| + \|D^{(i+1)} - D^{(i)}\| \\ &\leq (1 - \frac{1}{2^i}) \eta_2 + \frac{1}{2^{i+1}} \eta_2 \leq (1 - \frac{1}{2^{i+1}}) \eta_2 \\ \mu^{(i+1)} &\leq \mu^{(i)} + \frac{1}{2} \eta_2^{(i)} \leq \mu + (1 - \frac{1}{2^{i+1}}) \eta_2 \\ \sigma^{(i+1)} &\geq \sigma^{(i)} - \frac{1}{2} \eta_2^{(i)} \geq \frac{1}{2} (1 + \frac{1}{2^i} - \frac{1}{2^{i+1}}) \sigma \geq \frac{1}{2} (1 + \frac{1}{2^{i+1}}) \sigma \\ \eta_1^{(i+1)} &\leq \eta_1^{(i)} + \frac{1}{2} \eta_2^{(i)} \leq \eta_1 + (1 - \frac{1}{2^{i+1}}) \eta_2 \\ \eta_2^{(i+1)} &\leq \frac{1}{2} \eta_2^{(i)} \leq \frac{1}{2^{i+1}} \eta_2. \end{aligned}$$

This completes the induction.

Applying the induction to the sequence starting at  $D^{(i)}$ , we have for every  $j \geq 0$ ,

$$\|D^{(i+j)} - D^{(i)}\| \leq (1 - \frac{1}{2^{j+1}}) \eta_2^{(i)} \leq (1 - \frac{1}{2^{j+1}}) \frac{1}{2^i} \eta_2.$$



This shows that  $D^{(i)}$  is a Cauchy sequence that tends to a limit  $D^{(\infty)}$  with  $\|D^{(\infty)} - D\| \leq \eta_2$ . From this inequality, we also deduce that  $\|D^{(\infty)} - D^{(i)}\| \leq \frac{1}{2^{i+1}} \eta_2$ , so  $D^{(i)}$  converges geometrically to  $D^{(\infty)}$ .

Moreover, for each  $i$ , we have  $\varepsilon^{(i)} = \|E^{(i)}\| \leq \sigma^{-1} \eta_2^{(i)} \leq \frac{1}{2^i} \sigma^{-1} \eta_2$ . Hence, the matrix

$$\hat{E} = (1 + E^{(0)}) (1 + E^{(1)}) (1 + E^{(2)}) \dots - 1$$

is well defined, and

$$\begin{aligned} \log(1 + \|\hat{E}\|) &\leq \log(1 + \varepsilon^{(0)}) + \log(1 + \varepsilon^{(1)}) + \log(1 + \varepsilon^{(2)}) + \dots \\ &\leq 2 \sigma^{-1} \eta_2. \end{aligned}$$

We deduce that

$$\|\hat{E}\| \leq e^{2\sigma^{-1}\eta_2} - 1 \leq 3 \sigma^{-1} \eta_2,$$

since  $\sigma^{-1} \eta_2 \leq \frac{1}{32}$ .

We claim that  $M^{(i)} = D^{(i)} + H^{(i)}$  converges geometrically to

$$M^{(\infty)} = (1 + \hat{E})^{-1} M^{(0)} (1 + \hat{E}).$$

For any matrix  $M, E \in \mathbb{C}^{n \times n}$  with  $\|E\| < \varepsilon < 1$ , we have

$$\begin{aligned} \|(1 + E)^{-1} M (1 + E) - M\| &= \|ME - E(1 + E)^{-1} M (1 + E)\| \\ &\leq \|M\| (\varepsilon + \varepsilon (1 + \varepsilon) \|(1 + E)^{-1}\|) \\ &\leq \varepsilon \|M\| (1 + (1 + \varepsilon) (1 - \varepsilon)^{-1}) \\ &= \frac{2\varepsilon}{1 - \varepsilon} \|M\|. \end{aligned} \tag{7}$$

Let  $\hat{E}^{(i)} = (1 + E^{(i)}) (1 + E^{(i+1)}) (1 + E^{(i+2)}) \dots - 1$ . By the same arguments as above, we have  $\hat{\varepsilon}_i := \|E^{(i)}\| \leq 3 \sigma^{-1} \eta_2^{(i)} = \frac{3}{2^{i+1}} \sigma^{-1} \eta_2$ . Since  $M^{(\infty)} = (1 + \hat{E}^{(i)})^{-1} M^{(i)} (1 + \hat{E}^{(i)})$ , the inequality (7) implies

$$\begin{aligned} \|M^{(\infty)} - M^{(i)}\| &\leq \frac{2\hat{\varepsilon}_i}{1 - \hat{\varepsilon}_i} (\|D^{(i)}\| + \|H^{(i)}\|) \\ &\leq \frac{2\hat{\varepsilon}_i}{1 - \hat{\varepsilon}_i} (\mu_i + \eta_1^{(i)} + \eta_2^{(i)}) \\ &\leq \frac{3}{2^i} \frac{\sigma^{-1} \eta_2}{1 - \hat{\varepsilon}_i} (\mu + \eta_1 + \eta_2). \end{aligned}$$

This shows that  $M^{(i)}$  converges geometrically to  $M^{(\infty)}$ . We deduce that the sequence  $H^{(i)} = M^{(i)} - D^{(i)}$  also converges geometrically to a limit  $H^{(\infty)}$  with  $\|H^{(\infty)}\| \leq \eta_1 + \eta_2$ . Since  $\lim_{i \rightarrow \infty} \eta_2^{(i)} = 0$ , we finally observe that  $M^{(\infty)} = D^{(\infty)} + H^{(\infty)}$  is block diagonal.  $\square$

**THEOREM 4.** Assume  $I_k = \{k\}$  for all  $k$ . Then, under the assumptions of Theorem 3, the sequence  $(D, H), (D', H'), (D'', H''), \dots$  converges quadratically to  $(D^{(\infty)}, 0)$ .

**Proof.** The extra assumption implies that  $\eta_1^{(i)} = 0$  for all  $i$ . Let us show by induction over  $i$  that we now have

$$\eta_2^{(i)} \leq \frac{1}{2^{2^i - 1}} \eta_2.$$

This is clear for  $i = 0$ . Assume that the result holds for a given  $i$ . Then we may apply Lemma 2 to  $(D^{(i)}, H^{(i)})$  for  $\delta = 2^{-2^i+1}$ , and obtain

$$\begin{aligned}\eta_2^{(i+1)} &\leq \frac{1}{2^{2^i-1}} \eta_2^{(i)} \\ &\leq \frac{1}{2^{2^{i+1}-1}}.\end{aligned}$$

Since  $\|D^{(i+1)} - D^{(i)}\| \leq \eta_2^{(i)}$ , this establishes the quadratic convergence.  $\square$

## 4. ALGORITHMS

### 4.1. Clustering

Let  $M = D + H$  be the perturbation of a diagonal matrix (5) as in the previous section. In order to apply theorem 3, we first have to find a suitable clustering (4). For a given threshold separation  $\delta$ , we will simply take the finest clustering (i.e. for which  $p$  is maximal) with the property that  $|\lambda_i - \lambda_j| \leq \delta \Rightarrow i \sim j$ . This clustering can be computed using the algorithm Cluster below.

---

#### Algorithm Cluster

**Input:** eigenvalues  $\lambda_1, \dots, \lambda_n \in \mathbb{B}$  and  $\delta \in \mathbb{F}^{\geq}$

**Output:** the finest clustering (4) with  $[\lambda_i - \lambda_j] \leq \delta \Rightarrow i \sim j$

---

- Let  $G$  be the graph with vertices  $1, \dots, n$  and such that  $i$  and  $j$  are connected if and only if  $[\lambda_i - \lambda_j] \leq \delta$ .
  - Let  $H$  be the transitive closure of  $G$ .
  - Let  $H_1, \dots, H_p$  the connected components of  $H$ .
  - Let  $I_k$  be the set of vertices of  $H_k$  for each  $k$ .
- 

### 4.2. Certification in the case of perturbed diagonal matrices

In order to apply theorem 3, it now remains to find a suitable threshold  $\delta$  for which the conditions of the theorem hold. Starting with  $\delta = 0$ , we will simply increase  $\delta$  to  $\sigma(D)$  whenever the conditions are not satisfied. This will force the number  $p$  of clusters to decrease by at least one at every iteration, whence the algorithm terminates. Notice that the workload of one iteration is  $O(n^2)$ , so the total running time remains bounded by  $O(n^3)$ .

---

#### Algorithm DiagonalCertify

**Input:** a diagonal ball matrix  $D \in \mathbb{B}^{n \times n}$  with entries  $\lambda_1, \dots, \lambda_n$  and an off diagonal ball matrix  $H \in \mathbb{B}^{n \times n}$

**Output:** a clustering  $I$  and  $\hat{\varepsilon} \in \mathbb{F}$  such that, for any  $M \in D$  and  $H \in H$ , the conditions of theorem 3 hold and  $\|\hat{E}\| \leq \hat{\varepsilon}$

---

$\delta := 0$

Repeat

    Compute the clustering  $I$  for  $\lambda_1, \dots, \lambda_n$  and  $\delta$  using Cluster

    Let  $\mu := \|D\|$ ,  $\sigma := \sigma^I(D)$ ,  $\eta_1 := \|\Delta^I(H)\|$  and  $\eta_2 := \|\Omega^I(H)\|$

    Let  $\alpha := \min \left\{ \frac{\sigma}{6\mu}, \frac{1}{4} \right\}$

    If  $\lceil \eta_1 + \eta_2 \rceil \leq \lfloor \frac{\alpha\mu}{8} \rfloor$  and  $\lceil \eta_2 \rceil \leq \lfloor \frac{\alpha\sigma}{8} \rfloor$ , then return  $(I, \lceil \frac{3\eta_2}{\sigma} \rceil)$

    Set  $\delta := \lceil \sigma \rceil$

---

### 4.3. Certification of approximate eigenvectors and eigenvalues

Let us now return to the original problem of certifying a numerical solution to an eigenproblem. We will denote by  $\mathbb{1}_n$  the  $n \times n$  matrix of which all entries are one.

---

**Algorithm EigenvectorCertify**


---

**Input:**  $M = \mathcal{B}(M_c, M_r) \in \mathbb{B}^{n \times n}$  and  $T_c \in \mathbb{F}[\mathbf{i}]^{n \times n}$

such that  $T_c^{-1} M_c T_c$  is approximately diagonal

**Output:** a clustering  $I$  and  $T = \mathcal{B}(T_c, T_r) \in \mathbb{B}^{n \times n}$  such that for any  $M \in \mathbf{M}$ , there exists a  $T \in \mathbf{T}$  for which  $T^{-1} M T$  is block diagonal

---

Compute  $D := \mathcal{B}(T_c, 0)^{-1} M \mathcal{B}(T_c, T_r)$

Let  $(I, \varepsilon) := \text{DiagonalCertify}(\Delta^*(D), \Omega^*(D))$

Let  $E := \mathcal{B}(T_c, 0) \mathcal{B}(0, \varepsilon) \mathbb{1}_n$

Let  $(T_r)_{i,j} := \lceil E_{i,j} \rceil$  for all  $i, j$

Return  $(I, \mathcal{B}(T_c, T_r))$

---

Obviously, any eigenvalue  $\lambda \in \mathbb{C}$  of a matrix  $M \in \mathbb{C}^{n \times n}$  satisfies  $|\lambda| \leq \|M\|$ . We may thus use the following modification of EigenvectorCertify in order to compute enclosures for the eigenvalues of  $M$ .

---

**Algorithm EigenvalueCertify**


---

**Input:**  $M = \mathcal{B}(M_c, M_r) \in \mathbb{B}^{n \times n}$  and  $T_c \in \mathbb{F}[\mathbf{i}]^{n \times n}$

such that  $T_c^{-1} M_c T_c$  is approximately diagonal

**Output:** ball enclosures  $\lambda_1, \dots, \lambda_n \in \mathbb{B}$  for the eigenvalues of  $M$ , with the appropriate multiplicities in cases of overlapping

---

Compute  $D := \mathcal{B}(T_c, 0)^{-1} M \mathcal{B}(T_c, T_r)$

Let  $(I, \varepsilon) := \text{DiagonalCertify}(\Delta^*(D), \Omega^*(D))$

Let  $\eta_1 := \|\Delta^I(\Omega^*(D))\|$  and  $\eta_2 := \|\Omega^I(\Omega^*(D))\|$

For each  $k \in \{1, \dots, p\}$  do

    If  $I_k = \{i\}$  for some  $i$ , then let  $\lambda_i := \mathcal{B}((D_c)_{i,i}, \lceil \eta_2 \rceil)$

    Otherwise

        Let  $c$  be the barycenter of the  $D_{i,i}$  with  $i \in I_k$

        Let  $r$  be the maximum of  $|D_{i,i} - c|$  for  $i \in I_k$

        Let  $\lambda_i := c + \mathcal{B}(0, \lceil r + \eta_1 + 2\eta_2 \rceil)$  for all  $i \in I_k$

Return  $(\lambda_1, \dots, \lambda_n)$

---

## 5. POSSIBLE EXTENSIONS

Let  $M \in \mathbb{C}^{n \times n}$  be a matrix with a (numerically) multiple eigenvalue  $\lambda$ . We have already stressed that it is generally impossible to provide non trivial certifications for the corresponding eigenvectors. Nevertheless, two observations should be made:

- If the eigenspace  $E_\lambda$  corresponding to  $\lambda$  has dimension 1, then small perturbations of the matrix  $M$  only induce small perturbations of  $\lambda$  and  $E_\lambda$ .
- Let  $F_\lambda$  denote the full invariant subspace associated to the eigenvalue  $\lambda$  (or all eigenvalues in the cluster of  $\lambda$ ). Then small perturbations of  $M$  only induce small perturbations of  $\lambda$  and  $F_\lambda$ .

More precisely, in these two cases, we may search for ball enclosures for orthonormal bases of the vector spaces  $E_\lambda$  resp.  $F_\lambda$ , which do not contain the zero vector.

When considering the numeric solution (1) of the eigenproblem for  $M$ , the column vectors which generate  $F_\lambda$  are usually far from being orthogonal. Orthonormalization can only be done at the expense of making  $T^{-1} M T$  only upper triangular. Moreover, the orthogonalization implies a big loss of accuracy, which requires the application of a correction method for restoring the accuracy. It seems that the fundamental Newton iteration from Section 3.2 can actually be used as a correction method. For instance, for small perturbations of the matrix

$$D = \begin{pmatrix} \lambda_1 & 1 & 0 & 0 \\ 0 & \lambda_1 & 0 & 0 \\ 0 & 0 & \lambda_2 & 1 \\ 0 & 0 & 0 & \lambda_2 \end{pmatrix},$$

it can be shown that the fundamental iteration still converges. However, for more general block diagonal matrices with triangular blocks, the details are quite technical and yet to be worked out.

Yet another direction for future investigations concerns the quadratic convergence. As a refinement of Lemma 1, we might replace  $D$  by a block diagonal matrix with entries  $\Lambda_1, \dots, \Lambda_p$ . Instead of taking  $B_{i,j} = \frac{M_{i,j}}{\lambda_j - \lambda_i}$ , we then have to solve equations of the form

$$B_{i,j} \Lambda_j - \Lambda_i B_{i,j} = M_{i,j}.$$

If the  $\Lambda_i$  are sufficiently close to  $\lambda_i \text{Id}$ , it might then be possible to adapt the fundamental iteration accordingly so as to achieve quadratic convergence for the strongly off diagonal part.

## BIBLIOGRAPHY

- [1] G. Alefeld and J. Herzberger. *Introduction to interval analysis*. Academic Press, New York, 1983.
- [2] L. Blum, F. Cucker, M. Shub, and S. Smale. *Complexity and real computation*. Springer-Verlag, 1998.
- [3] J. J. Dongarra, C. B. Moler, and J. H. Wilkinson. Improving the accuracy of computed eigenvalues and eigenvectors. *SIAM Journal on Numerical Analysis*, 20(1):23–45, 1983.
- [4] G. H. Golub and F. Van Loan. *Matrix Computations*. JHU Press, 1996.
- [5] S. Graillat and Ph. Trébuchet. A new algorithm for computing certified numerical approximations of the roots of a zero-dimensional system. In *Proc. ISSAC '09*, pages 167–174. ACM Press, 2009.
- [6] J. van der Hoeven. Ball arithmetic. Technical report, HAL, 2009. <http://hal.archives-ouvertes.fr/hal-00432152>.
- [7] L. Jaulin, M. Kieffer, O. Didrit, and E. Walter. *Applied interval analysis*. Springer, London, 2001.
- [8] F. Johansson. Arb: A c library for ball arithmetic. *ACM Commun. Comput. Algebra*, 47(3/4):166–169, 2014.
- [9] U. W. Kulisch. *Computer Arithmetic and Validity. Theory, Implementation, and Applications*. Number 33 in Studies in Mathematics. de Gruyter, 2008.
- [10] S. Miyajima. Fast enclosure for all eigenvalues in generalized eigenvalue problems. *Journal of Computational and Applied Mathematics*, 233(11):2994–3004, 2010.
- [11] R. E. Moore. *Interval Analysis*. Prentice Hall, Englewood Cliffs, N.J., 1966.
- [12] A. Neumaier. *Interval methods for systems of equations*. Cambridge university press, Cambridge, 1990.
- [13] S. Rump and S. Graillat. Verified error bounds for multiple roots of systems of nonlinear equations. *Num. Algs.*, 54:359–377, 2010.
- [14] S. M. Rump. Guaranteed inclusions for the complex generalized eigenproblem. *Computing*, 42(2):225–238, 1989.
- [15] S. M. Rump. INTLAB - INTerval LABoratory. In Tibor Csendes, editor, *Developments in Reliable Computing*, pages 77–104. Kluwer Academic Publishers, Dordrecht, 1999. <http://www.ti3.tu-harburg.de/rump/>.
- [16] S. M. Rump. Computational error bounds for multiple or nearly multiple eigenvalues. *Linear algebra and its applications*, 324(1-3):209–226, 2001.
- [17] T. Yamamoto. Error bounds for computed eigenvalues and eigenvectors. *Numerische Mathematik*, 34(2):189–199, 1980.